
Optimizing Compilers

Inter-Procedural Dataflow Analysis

Markus Schordan

Institut für Computersprachen
Technische Universität Wien

Syntax

$$P_{\star} ::= \text{begin } D_{\star} \ S_{\star} \ \text{end}$$
$$D ::= D; D \mid \text{proc } p(\text{val } x; \text{res } y) \text{ is}^{\ell_n} S \text{ end}^{\ell_x}$$
$$S ::= \dots \mid [\text{call } p(a, z)]_{\ell_r}^{\ell_c}$$

Labeling scheme

- procedure declarations
 - ℓ_n : for entering the body
 - ℓ_x : for exiting the body
- procedure calls
 - ℓ_c : for the call
 - ℓ_r : for the return

Analysing Procedures

We consider procedures with call-by-value and call-by-result parameters.

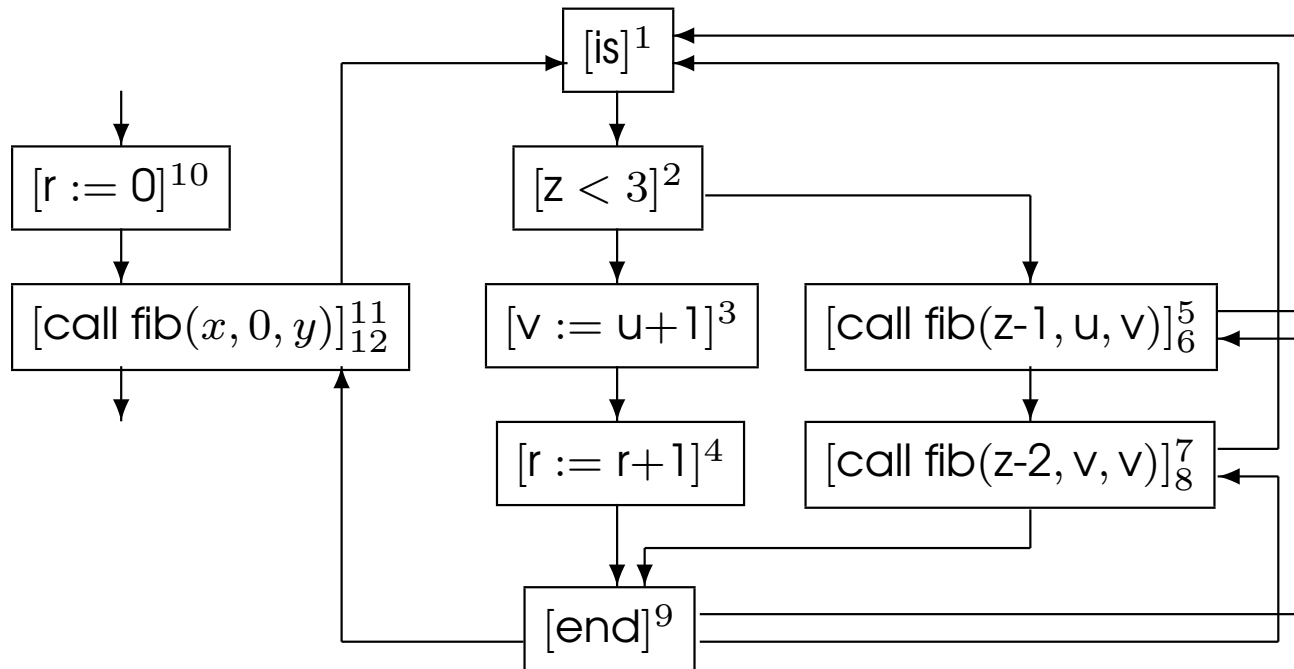
Example:

```
begin
  proc fib(val z,u; res v) is
    if z<3 then
      (v:=u+1; r:=r+1)
    else (
      call fib (z-1,u,v);
      call fib (z-2,v,v)
    )
  end;
  r:=0;
  call fib(x,0,y)
end
```

Example Flow Graph

main

proc fib(val z, u; res v)



Flow Graph for Procedures

	$[\text{call } p(a, z)]_{l_r}^{l_c}$	$\text{proc } p(\text{val } x; \text{res } y) \text{ is}^{l_n} S \text{ end}^{l_x}$
init	l_c	l_n
final	$\{l_r\}$	$\{l_x\}$
blocks	$\{[\text{call } p(a, z)]_{l_r}^{l_c}\}$	$\{\text{is}^{l_n}\} \cup \text{blocks}(S) \cup \{\text{end}^{l_x}\}$
labels	$\{l_c, l_r\}$	$\{l_c, l_r\} \cup \text{labels}(S)$
flow	$\{(l_c; l_n), (l_x; l_r)\}$	$\{(l_n, \text{init}(S))\} \cup \text{flow}(S) \cup \{l, l_x \mid l \in \text{final}(S)\}$

- $(l_c; l_n)$ is the flow corresponding to **calling** a procedure at l_c and entering the procedure body at l_n and
- $(l_x; l_r)$ is the flow corresponding to exiting a procedure body at l_x and **returning** to the call at l_r .

Naive Formulation

Treat the three kinds of flow, (ℓ_1, ℓ_2) , $(\ell_c; \ell_n)$, $(\ell_x; \ell_r)$ in the same way.

Equation system:

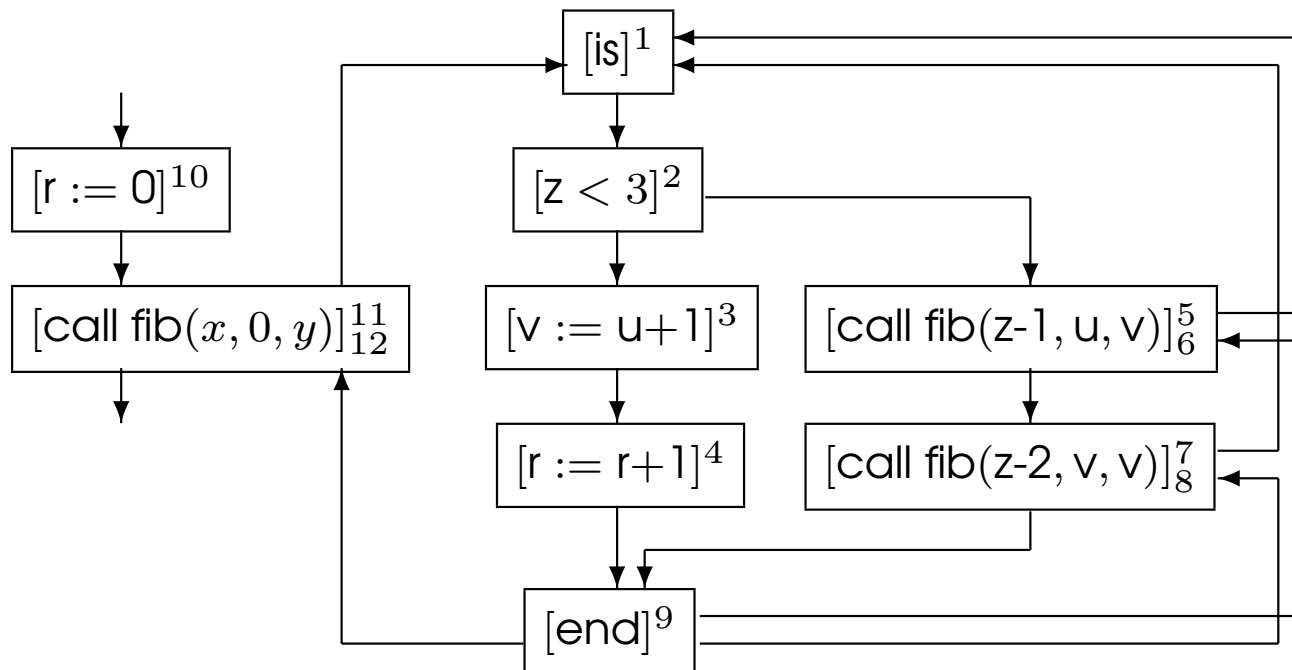
$$\begin{aligned} A_o(\ell) &= \sqcup \{A_\bullet(\ell') \mid (\ell', \ell) \in F \vee (\ell'; \ell) \in F\} \sqcup \iota_E^\ell \\ A_\bullet(\ell) &= f_\ell^A(A_o(\ell)) \end{aligned}$$

- both procedure calls $(\ell_c; \ell_n)$ and procedure returns $(\ell_x; \ell_r)$ are treated like “goto’s”.
- there is no mechanism for ensuring that information flowing along $(\ell_c; \ell_n)$ flows back along $(\ell_x; \ell_r)$ to the *same* call
- intuitively, the equation system considers a much too large set of “paths” through the program and hence will be grossly imprecise (although formally on the safe side)

Matching Procedure Entries and Exits

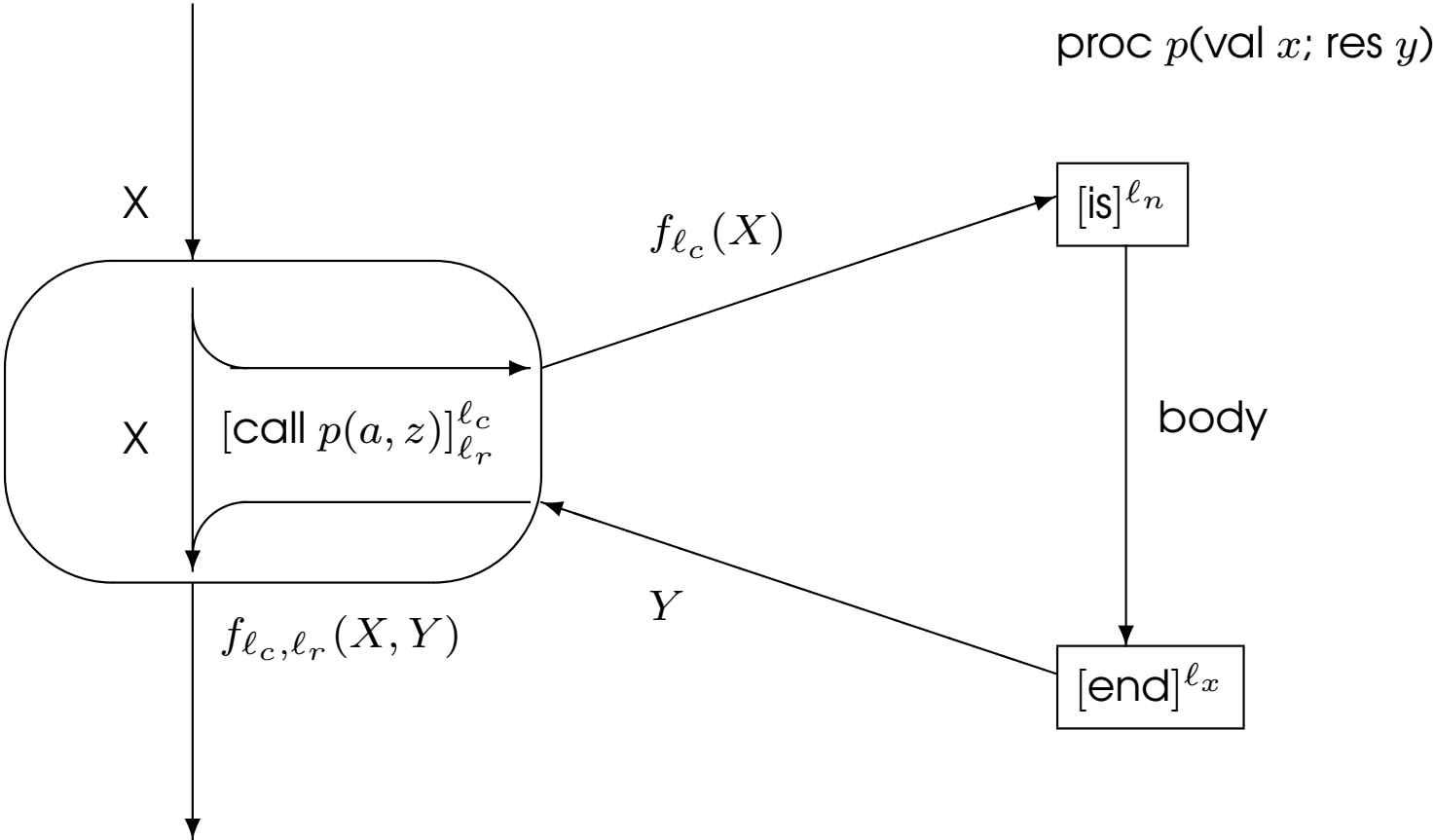
main

proc fib(val z, u; res v)



We want to overcome the shortcoming of the naive formulation by restricting attention to paths that have the proper nesting of procedure calls and exits.

General Formulation: Calls and Returns



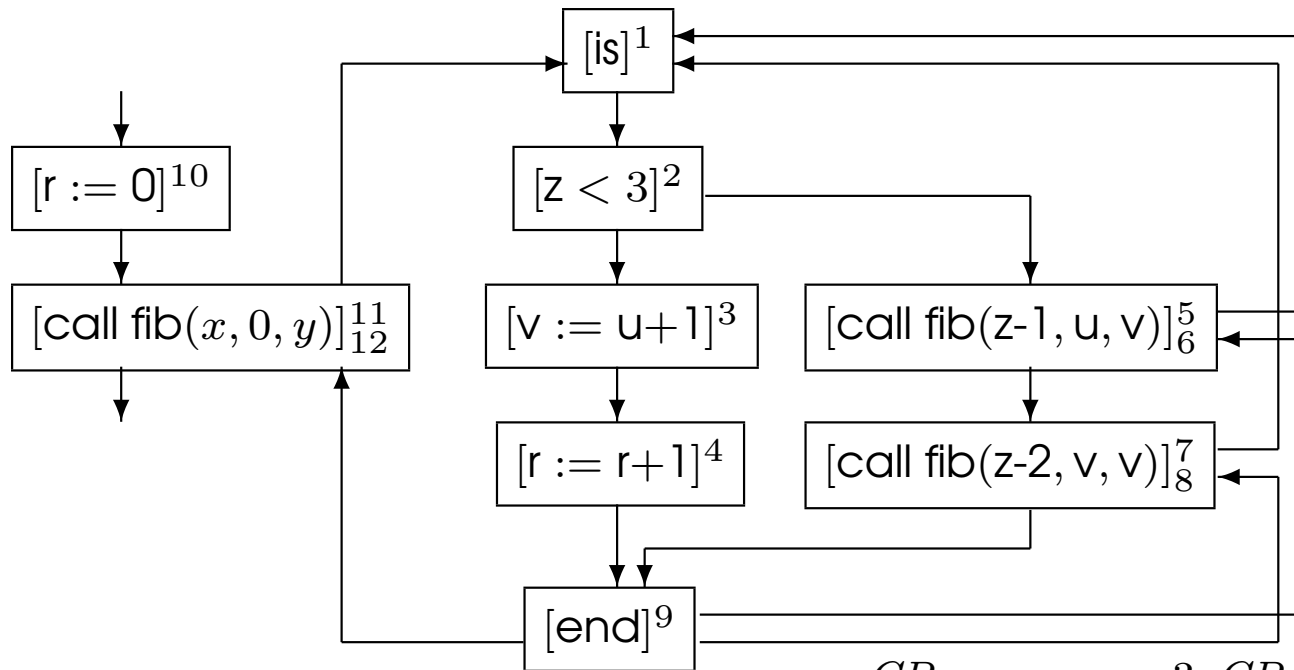
“Meet” over Valid Paths (MVP)

A complete path from l_1 to l_2 in P_\star has proper nesting of procedure entries and exits; and a procedure returns to the point where it was called:

$$\begin{array}{ll} CP_{l_1, l_2} \longrightarrow l_1 & \text{whenever } l_1 = l_2 \\ CP_{l_1, l_3} \longrightarrow l_1, CP_{l_2, l_3} & \text{whenever } (l_1, l_2) \in \text{flow}_\star \\ CP_{l_c, l} \longrightarrow l_c, CP_{l_n, l_x}, CP_{l_r, l} & \text{whenever } P_\star \text{ contains } [\text{call } p(a, z)]_{l_r}^{l_c} \\ & \text{and } \text{proc } p(\text{val } x; \text{res } y) \text{ is }^{l_n} S \text{ end}^{l_x} \end{array}$$

Definition: $(l_c, l_n, l_r, l_x) \in \text{interflow}_\star$ if P_\star contains $[\text{call } p(a, z)]_{l_r}^{l_c}$ as well as $\text{proc } p(\text{val } x; \text{res } y) \text{ is }^{l_n} S \text{ end}^{l_x}$

Example



$CP_{10,12} \rightarrow 10, CP_{11,12}$
 $CP_{11,12} \rightarrow 11, CP_{1,9}, CP_{12,12}$
 $CP_{1,9} \rightarrow 1, CP_{2,9}$
 $CP_{2,9} \rightarrow 2, CP_{3,9}$
 $CP_{2,9} \rightarrow 2, CP_{5,9}$

$CP_{3,9} \rightarrow 3, CP_{4,9}$
 $CP_{4,9} \rightarrow 4, CP_{9,9}$
 $CP_{5,9} \rightarrow 5, CP_{1,9}, CP_{6,9}$
 $CP_{6,9} \rightarrow 6, CP_{7,9}$
 $CP_{7,9} \rightarrow 7, CP_{1,9}, CP_{8,9}$
 $CP_{8,9} \rightarrow 8, CP_{9,9}$

Some valid paths: (10,11,1,2,3,4,9,12) and (10,11,1,2,5,1,2,3,4,9,6,7,1,2,3,4,9,8,9,12)

A non-valid path: (10,11,1,2,5,1,2,3,4,9,12)

Valid Paths

A **valid path** starts at the entry node init_* of P_* , all the procedure exits match the procedure entries but some procedures might be entered but not yet exited:

$$VP_* \longrightarrow VP_{\text{init}_*, l}$$

whenever $l \in \text{Lab}_*$

$$VP_{l_1, l_2} \longrightarrow l_1$$

whenever $l_1 = l_2$

$$VP_{l_1, l_3} \longrightarrow l_1, VP_{l_2, l_3}$$

whenever $(l_1, l_2) \in \text{flow}_*$

$$VP_{l_c, l} \longrightarrow l_c, CP_{l_n, l_x}, VP_{l_r, l}$$

whenever P_* contains $[\text{call } p(a, z)]_{l_r}^{l_c}$
and $\text{proc } p(\text{val } x; \text{res } y) \text{ is }^{l_n} S \text{ end}^{l_x}$

$$VP_{l_c, l} \longrightarrow l_c, VP_{l_n, l}$$

whenever P_* contains $[\text{call } p(a, z)]_{l_r}^{l_c}$
and $\text{proc } p(\text{val } x; \text{res } y) \text{ is }^{l_n} S \text{ end}^{l_x}$

MVP Solution

$$MVP_{\circ}(\ell) = \bigsqcup \{f_{\vec{\ell}}(\iota) \mid \vec{\ell} \in vpath_{\circ}(\ell)\}$$

$$MVP_{\bullet}(\ell) = \bigsqcup \{f_{\vec{\ell}}(\iota) \mid \vec{\ell} \in vpath_{\bullet}(\ell)\}$$

where

$$vpath_{\circ}(\ell) = \{[\ell_1, \dots, \ell_{n-1}] \mid n \geq 1 \wedge \ell_n = \ell \wedge [\ell_1, \dots, \ell_n] \text{ is valid path}\}$$

$$vpath_{\bullet}(\ell) = \{[\ell_1, \dots, \ell_n] \mid n \geq 1 \wedge \ell_n = \ell \wedge [\ell_1, \dots, \ell_n] \text{ is valid path}\}$$

The MVP solution may be undecidable for lattices satisfying the Ascending Chain Condition, just as was the case for the MOP solution.

Making Context Explicit

- The MVP solution may be undecidable for lattices of finite height (as was the case for the MOP solution)
- We have to reconsider the MFP solution and avoid taking too many invalid paths into account
- Encode information about the paths taken into data flow properties themselves
- Introduce context information

MFP Counterpart

Context sensitive analysis: add context information

- call strings:
 - an abstraction of the sequences of procedure calls that have been performed so far
 - example: the program point where the call was initiated
- assumption sets:
 - an abstraction of the states in which previous calls have been performed
 - example: an abstraction of the actual parameters of the call

Context insensitive analysis: take no context information into account.

Call Strings as Context

- Encode the path taken
- Only record flows of the form (ℓ_c, ℓ_n) corresponding to a procedure call
- we take as context $\Delta = \text{Lab}^*$ where the most recent label ℓ_c of a procedure call is at the right end
- Elements of Δ are called **call strings**
- The sequence of labels $\ell_c^1, \ell_c^2, \dots, \ell_c^n$ is the call string leading to the current call which happened at ℓ_c^1 ; the previous calls where at $\ell_c^2 \dots \ell_c^n$. If $n = 0$ then no calls have been performed so far.

For the example program the following call strings are of interest:

$\Lambda, [11], [11, 5], [11, 7], [11, 5, 5], [11, 5, 7], [11, 7, 5], [11, 7, 7], \dots$

Abstracting Call Strings

Problem: call strings can be arbitrarily long (recursive calls)

Solution: truncate the call strings to have length of at most k for some fixed number k

- $\Delta = \text{Lab}^{\leq k}$
- $k = 0$: context insensitive analysis
 - Λ (the call string is the empty string)
- $k = 1$: remember the last procedure call
 - $\Lambda, [11], [5], [7]$
- $k = 2$: remember the last two procedure calls
 - $\Lambda, [11], [11, 5], [11, 7], [5, 5], [5, 7], [7, 5], [7, 7]$

References

- Material for this 4th lecture (part 2)
www.complang.tuwien.ac.at/markus/optub.html
- Book
Flemming Nielson, Hanne Riis Nielson, Chris Hankin:
Principles of Program Analysis.
Springer, (450 pages, ISBN 3-540-65410-0), 1999.
– Chapter 2 (Data Flow Analysis)